

# Machine Learning Project Document

## Multi-Class Classification: Predicting University Student Grades

### 1. Problem Statement

This project aims to develop a system that predicts the final grade of university students as one of several categories (A, B, C, D, F). The prediction is based on academic and behavioral indicators such as attendance, participation, and assessment scores. The model should help identify students at academic risk early.

### 2. Dataset Description

Source: University internal records or synthetic data

Format: CSV / Excel

### 3. Data Preprocessing

This step ensures the data is clean, consistent, and ready for machine learning. It involves treating missing values, converting categorical data to numeric, normalizing feature values, and preparing the data splits for training and testing.

### 4. Exploratory Data Analysis (EDA)

EDA is performed to understand the data distribution, identify patterns, and detect anomalies or outliers. It provides insights into feature importance, correlation, and class balance, guiding model selection and refinement.

### 5. Model Selection & Justification

This step involves choosing appropriate algorithms for multi-class classification. The decision is based on the problem characteristics, data size, and performance requirements. The goal is to select a model that balances accuracy, interpretability, and robustness.

### 6. Model Training & Tuning

The chosen model(s) are trained on the prepared dataset, and their hyperparameters are optimized for better performance. Metrics such as accuracy, precision, recall, and F1-score are used to evaluate the effectiveness of each model during training.

### 7. Test Set Evaluation

The final model is evaluated on a separate test set to assess its generalization ability. This evaluation helps understand how well the model will perform on unseen data. Performance is measured using relevant metrics for multi-class classification.

## 8. Limitations and Future Work

This section identifies the known limitations of the current model, such as data imbalance or missing contextual factors like course difficulty. It also proposes improvements, such as adding more features or refining the model with real-time data.

## 9. File Structure

A well-organized project directory should include folders for raw data, cleaned data, scripts, model artifacts, notebooks, and documentation. This structure ensures reproducibility and clarity.

## 10. Model Reusability and Testing Requirement

The system must support the ability to save trained models, load them for future predictions, and test them on new data. This ensures that models are reusable, and predictions can be made consistently outside the training phase.

## 11. GitHub Project Requirement

The complete project must be uploaded to a GitHub repository. It should include:

- A clear and descriptive README file
- Structured directories for data, code, and models
- Instructions for running the project and making predictions
- License and contribution guidelines (if applicable)

## Unsupervised Learning Alternative: Student Performance Clustering

In addition to the supervised classification task, this project can also be approached as an unsupervised learning problem. The goal is to group students into clusters based on similarities in their academic behavior and performance, without using predefined labels such as final grades.

### 1. Objective

Discover natural groupings of students (e.g., high performers, at-risk students, average performers) based on features

### 2. Preprocessing

Clean and normalize the dataset as done for the supervised version. Ensure features are on a similar scale for effective clustering.

### 3. Clustering Algorithm Selection

Choose appropriate clustering algorithms based on the data shape and expected distribution.

### 4. Cluster Evaluation

Use unsupervised metrics to evaluate clustering quality. Also perform a qualitative analysis by interpreting the cluster centers and comparing student characteristics.

## **5. Interpretation and Application**

Analyze the resulting clusters to identify trends and behaviors

These clusters can help inform academic support strategies and future interventions.